# Experimental Analysis of a Ring Oscillator Network for Hardware Trojan Detection in a 90nm ASIC

Andrew Ferraiuolo, Xuehui Zhang, and Mohammad Tehranipoor
ECE, University of Connecticut
{andrew.ferraiuolo,xuehui.zhang,tehrani}@engr.uconn.edu

## ABSTRACT

The modern integrated circuit (IC) manufacturing process has exposed chip designers to hardware Trojans which threaten circuits bound for critical applications. This paper details the implementation and analysis of a novel ring oscillator network technique for Trojan detection in an application specific integrated circuit (ASIC). The ring oscillator network serves as a power supply monitor by detecting fluctuations in characteristic frequencies due to malicious modifications (i.e. hardware Trojans) in the circuit under authentication. The ring oscillator network was implemented and fabricated in 40 IBM 90nm ASICs with controlled hardware Trojans. This work analyzes the impact of Trojans with varied partial activity, area, and location on the proposed ring oscillator structure and demonstrates that stealthy Trojans can be efficiently detected with this technique even while obfuscated by process variations, background noise, and environment noise.

## Categories and Subject Descriptors

B.8.0 [**Performance and Reliability**]: General

## General Terms

Security

## Keywords

Hardware Trojan detection, IC trust, process variations, and on-chip measurement

## 1. INTRODUCTION

Recent changes to the integrated circuit (IC) manufacturing process and rapid globalization have made integrated circuit designers increasingly vulnerable to malicious modifications (i.e. hardware Trojans) [1]. A hardware Trojan may be implemented by the addition or omission of gates or by the modification of design parameters and may act to destroy or disable the chip, reduce the performance of the design, or leak confidential information among other possibilities. A taxonomy exists to categorize Trojans based on its

malicious effect, triggering mechanism, the abstraction level of the design, and the physical characteristics of the Trojan among other considerations [3].

IC designers lacking a foundry are vulnerable to a class of hardware Trojan attacks during which an adversary inserts a hardware Trojan at the untrusted fabrication facility. The discovery and prevention of this class of Trojan attacks is referred to as the *IC trust problem* [10]. Following fabrication, an IC undergoes functional and structural tests during which automatic test pattern generation (ATPG) produces a sequence of input combinations which form a subset of all possible input combinations. While Trojans may depend on the inputs and intermediate signals of the original design, an exhaustive test pattern is infeasible for an IC of modest size.

Therefore, it is unlikely that a hardware Trojan will be *fully activated*, launching its malicious payload and observably modifying the IC's behavior, from these normal testing procedures [10]. Structural tests, such as those based on stuck-at or bridging fault models, are based solely on the original untampered netlist, and thus cannot guarantee Trojan detection. Furthermore, it is possible to construct a Trojan which cannot be activated through any test pattern based on the original circuit (e.g. a Trojan which depends on temperature or a wireless receiver for activation).

However, test patterns can provoke *partial activation* during which some of a Trojan's gates transition consuming power and altering gate delays. Numerous detection techniques have leveraged these Trojan-induced changes to the IC's side-channel information eliminating the need for an exhaustive test by monitoring changes in transient power [4][5][6][10], current[7], and delay [6][8][9]. Notably, many of these techniques require a *golden IC* signature constructed from verified Trojan-free circuits and thus assume that it is possible to obtain such a signature (e.g. from destructive reverse-engineering performed after a set of side-channel measurements).

The problem is exacerbated by process variations, measurement noise, and environmental variations which also alter these side-channels, and thus, obfuscate Trojans and complicate detection. Techniques which aim to improve the chance of activating a Trojan have been proposed in [11][12][13][14]. These techniques are at a disadvantage when attempting to detect Trojans with very specific, rare conditions and are only capable of detecting the functional category of Trojans described in [15]. Since many of these techniques also improve the partial activity of a hardware Trojan, a composite technique which increases activity and simultaneously measures side-channels may be desired [12][14].

The on-chip ring oscillator network (RON) structure was proposed to detect hardware Trojans by utilizing ring oscillators (ROs) as sensors for power network noise [2]. The frequency of an RO is dependent on the power supply, thus by measuring changes in the frequency the malicious addition or omission of gates may be

Figure 1: The RON structure and topology [2].



Figure 2: Layout for the test chip design.

detected. Ring oscillators are inherently tamper resistant since an RO's frequency will vary across multiple measurements unlike a lookup table or a simple constant. This technique may be coupled with any other previously proposed technique [11][12][13][14] for improved detection.

In this work, the RON structure is analyzed with silicon results from 90nm integrated circuits which include the ISCAS'89 s9234 benchmark circuit to provide background activity that contributes to obfuscation. This paper demonstrates that RON is an effective technique for detecting stealthy hardware Trojans under process and environmental variations and analyzes the impact of variations in Trojan size and activity levels. The effect of Trojan location in relation to the ROs and the IC power distribution network and the ability of RON to determine the location of an attack is explored. Lastly, in addition to comparing the Trojan-inserted ICs against a signature, this work analyzes a scheme for classifying Trojan-inserted and Trojan-free chips and performs a false-positive analysis.

The rest of this paper is organized as follows: Section 2 describes the RON architecture and topology. Section 3 describes the design and implementation of the 90nm ICs under experimentation. Section 4 describes the experimental and data collection procedures. Section 5 provides extensive analyses and results on the ICs. Finally, concluding remarks are given in Section 6.

## 2. BACKGROUND: RON THEORY, ARCHITECTURE, AND TOPOLOGY

The frequency of an $n$-stage RO is:

$$f = \frac{\mu_g \times (V_{DD} - V_{TH})^\alpha}{2n \times k_g} \qquad (1)$$

where $\alpha$ is the velocity saturation index, $V_{DD}$ is the supply voltage, $V_{TH}$ is the threshold voltage, $\mu_g$ is the carrier mobility, and $k_g$ is a gate-dependant constant [2]. However, in the presence of a Trojan, the load is increased and an additional voltage drop $\Delta V_{TROJ}$, is introduced changing the frequency of the RO to:

$$f = \frac{\mu_g \times (V_{DD} - \Delta V_{TROJ} - V_{TH})^\alpha}{2n \times k_g} \qquad (2)$$

Therefore, changes in the frequency of a RO may be measured to detect the presence of a Trojan [2].

The RO network structure shown in Figure 1 contains several ROs to be used as power supply noise sensors. The number of ROs, $N_{RO}$, to be used may be adjusted based on the area of the chip, the power structure of the chip, and the area that may be used to implement the RON structure [2]. Each RO consists of $n-1$ inverters and 1 NAND gate to allow it to be enabled/disabled as needed. In addition to the ROs, a decoder and multiplexer are used to control which RO is enabled and which RO is sent to the counter, respectively. The output of the multiplexer is routed to a counter which determines the total number of oscillations over a number of clock cycles which is controlled by a timer (labeled cycle count). The frequency may then be determined from the oscillation count. In order to provoke Trojan partial activation (which we stress differs from full activation), a linear feedback shift register (LFSR) is used to supply random test patterns to the circuit while the frequency measurement is in progress. It is crucial that the same test patterns are used for each RO and in each chip under authentication.

The stages of each RO are to be placed vertically such that each stage is adjacent to a different standard cell. This topology intends to maximize the sensors' coverage of the power-supply network, and thus the sensors' sensitivity to the small noise produced by Trojans. Since ROs are composed primarily of loosely distributed inverters, the overhead of the ROs is anticipated to be very low. The area overheads of the decoder and multiplexer both scale logarithmically with the number of ROs whereas the size of the counter is dependent on the anticipated maximum RO frequency and the duration of the test period. Lastly, the area overhead of the LFSR is dependent on the total number of inputs to the circuit. However, modern ICs usually include an LFSR for built-in self-test procedures, in which case an additional LFSR is not needed for the RON. If an LFSR is not already present, an LFSR of fewer bits than the number of total inputs may broadcast each LFSR output to several inputs. Since the RON structure is enabled only during a test process and disabled for the lifetime of the chip, the power consumption during normal operation is negligible.

## 3. IC DESIGN AND IMPLEMENTATION

### 3.1 Test Chip Design

In order to analyze the effectiveness of the RON structure, 40 test chips were designed and fabricated using IBM 90nm technology through MOSIS. All chips used in this work were fabricated on the same wafer. The RON architecture is inserted into the ISCAS s9234 benchmark which represents the design to be protected in the test chip. Figure 2 shows the layout of the test chips with the RON structure composed of $N_{ro} = 8$ $n = 61$-stage ROs ($RO_j$ where $1 \leq j \leq 8$) with one NAND gate and 60 inverters each distributed across the chip. It is important to note that the areas labeled $RO_1$

Figure 3: Design of a hardware Trojan stage $T_i$.

Table 2: Estimation of Trojan area overheads and noise.

| Trojan Number | Transistors | Percent Area | Trojan to Background Circuit Switching Ratio |
|---|---|---|---|
| T1 | 26 | 0.12% | 0.11% |
| T2 | 52 | 0.23% | 0.22% |
| T3 | 78 | 0.35% | 0.33% |
| T4 | 104 | 0.47% | 0.45% |
| T5 | 130 | 0.58% | 0.56% |
| T6 | 156 | 0.70% | 0.67% |
| T7 | 182 | 0.81% | 0.78% |

to $RO_8$ show the broad area in which that RO is confined rather than the total area occupied by that RO. Ring oscillator stages are placed in each standard cell row in an intentionally, loosely distributed fashion that improves its coverage of the power distribution network. Therefore, these areas are also occupied by background circuit and control structure components and the area overhead of the oscillators is substantially lower than the labeled areas. The approximate locations of the seven Trojan stages ($T_i$ where $1 \leq i \leq 7$) are labeled as well. The number of RO stages was selected so that the maximum observed frequency would not exceed the $400MHz$ operating frequency of the 90nm counters used in this design. The distance between the two adjacent RO components is limited to 10 times of the width of the flip-flops. Based on this design rule and the area of the chip, 8 ROs were used.

The feedback polynomial of the LFSR used in our test chip is

$$X^7 + X^3 + 1 \qquad (3)$$

To conserve area, this design uses an LFSR with only 8-bits to generate patterns for the 36 input s9234 benchmark. A broadcasting technique is used to assign this 8-bit output to the 36 inputs. An 8-bit decoder and 8-bit multiplexer are used for RO selection.

A 16-bit counter is used to measure the number of oscillations observed in the test duration which is controlled by a timer. In this design, the test duration of 500 clock cycles was selected based on the technology node and test area overhead.

## 3.2 Hardware Trojan Design

Each IC contains seven combinational hardware Trojan designs which may be completely deactivated. Since this design is implemented in 90nm CMOS technology, the static power dissipation, and thus side-channel contribution is negligible when the Trojans are deactivated. By using a single-IC multiple-Trojan design we are able to not only carry out a more extensive set of Trojan impact tests, but we are also able to isolate the effect of process variations from the effect of inserted Trojans on RO characteristic frequencies. Further, since the static power is present in the Trojan-free case, it is neglected in comparisons to Trojan-inserted cases, and the detection results provide a lower-bound.

The gate-level implementation of a Trojan stage is shown in Figure 3 where $troout[i]$ is the output of the $i^{th}$ Trojan stage, $troout[i-1]$ is the output of the previous Trojan stage, and $troen[i]$ is the enable signal for the $i^{th}$ stage which also asserts all prior enable signals when enabled. Trojan $T_i$ contains $i$ stages consisting of $i \times (4AND + 1INV)$ gates where each stage $i-1$ is also enabled if stage $i$ is enabled. The first Trojan, $T_1$ is driven by the $200MHz$ clock signal at the location of signal $troout[0]$. Note that the Trojan, $T_i$, is not derived of the trigger-payload Trojan design used in [4][12][13]. Here, each Trojan gate transitions once per clock cycle, therefore, the partial activity of each of these Trojans is simply $5i$ partial activations per clock cycle. The average ratio of Trojan partial activation to background circuit activity is estimated in the fourth column of Table 2.

The s9234 benchmark consists of 211 D flip-flops, 3570 inverters, and 2027 other gates. The number of transistors used in the

s9234 benchmark is estimated in Table 1 by assuming each flip-flop consists of 8 NAND or NOR gates and 2 inverters. As mentioned earlier, there are a total of seven Trojans ($T_1$ to $T_7$) in this design. The area overhead of each Trojan is summarized in Table 2.

## 4. EXPERIMENTAL SETUP

During data collection, the IC is mounted on and wired to a prototyping board which includes a high-density serial connector. The serial connector allows the prototyping board to interface with a Xilinx Spartan-6 FPGA on a Digilent Nexys 3 board. The FPGA is programmed to control the test sequence supplied to the IC and transmit the outputs of the IC to a computer using an on-board USB-UART module. The complete setup is shown in Figure 4.

The nominal supply voltage of the pins of the IC is 2.5V. This is converted internally to the nominal core voltage of 1.2V using a voltage divider. Since the s9234 benchmark circuit used in this design is small compared to a modern IC, in order to emulate the tight power design of a modern circuit, an external voltage divider is used to supply the IC with 1.875V and the core with 0.9V which is greater than the 0.80V minimum core voltage. In practice, reducing the power supply voltage will reduce the background circuit switching activity and improve Trojan detection rates. Therefore, it is desirable to reduce the supply voltage during measurement.

The FPGA includes a state machine which sequences through each ring oscillator, begins a data collection trial, selects each 4-bit window of the counter output for the current ring oscillator, and transmits each 4-bit window as a hex digit over the USB-UART connection. The process is repeated for 10 trials on each ring oscillator of each IC. The IC is supplied with 1.875V using a voltage divider and the board's 2.5V peripheral power supply over the serial connection along with a 200MHz clock signal. Each trial lasts 500 clock cycles.

As discussed in Section 3, each of the 40 ICs contains $N_T = 7$ pre-inserted hardware Trojan designs. During Trojan-free data collection each hardware Trojan circuit is disabled, as is any Trojan not being analyzed. Since the designs are implemented with CMOS circuits, the static dissipation is negligibly low. Furthermore, since all Trojan measurements are compared to the Trojan-free results (which include static dissipation) the presented detection results provide a conservative lower bound.

Figure 4: Data collection setup including a Spartan 6 FPGA connected to a prototyping board through a serial connector. The chip under authentication is placed on the prototyping board.

Table 3: Summary of validation data

| Measurement Noise | 0.23% |
|---|---|
| Intradie Variation | 8.05% |
| Interdie Variation | 16.67% |
| Mean RO Frequency | 291MHz |

## 5. EXPERIMENTAL RESULTS AND ANALYSIS

The frequency of a single ring oscillator on a single IC was measured 10 times. The measurement noise is then calculated with

$$\frac{Max\{f_{Trial1},...,f_{Trial10}\} - Min\{f_{Trial1},...,f_{Trial10}\}}{0.1\sum_{m=1}^{10} f_{Trialm}} \quad (4)$$

for a single IC and a single ring oscillator where $f_{Trialm}$ is the $m^{th}$ repeated measurement of frequency for that RO. This is repeated for all ICs and all ROs and averaged resulting in a measurement noise of 0.23%.

The impact of intra-die variation on an RO's frequencies was analyzed by comparing a single RO on an IC with other ROs on that same IC. For a single IC, intra-die variation is calculated with

$$\frac{Max\{f_{RO_1},...,f_{RO_8}\} - Min\{f_{RO_1},...,f_{RO_8}\}}{0.125\sum_{j=1}^{8} f_{RO_j}} \quad (5)$$

where $f_{RO_j}$ is the frequency of the $j^{th}$ RO. This calculation is repeated for all ICs and averaged resulting in a mean intra-die variation impact on frequency of 8.05%.

Of the 40 fabricated ICs, 38 functioned correctly and the remaining faulty ICs are omitted. The impact of inter-die variation on the frequency of a ring oscillator was determined by selecting a single RO and comparing the frequency of this RO across each IC. For a single RO the inter-die variation is calculated with

$$\frac{Max\{f_{IC1},...,f_{IC38}\} - Min\{IC_1,...,f_{IC38}\}}{(1/38)\sum_{k=1}^{38} f_{ICk}} \quad (6)$$

where $f_{ICk}$ is the frequency of the individual RO of interest on the $k^{th}$ integrated circuit. This calculation is repeated for all ROs and averaged resulting in a mean inter-die variation impact on frequency of 16.67%. The average RO frequency of all ROs on all ICs was 291$MHz$. The maximum recorded frequency was 315$MHz$ which was less than the 400$MHz$ frequency the counter was timing closed at. These results are summarized in Table 3.

### 5.1 Trojan Impact Analysis



Figure 5: The impact of inserted hardware Trojans on RO frequencies isolated from process variations.

The direct impact of hardware Trojan induced power supply noise on ring oscillator frequencies is analyzed by measuring the frequency of each RO on each IC for the Trojan-free case as well as for each Trojan. The mean impact of a particular Trojan on a particular RO is then computed by comparing the frequency of that RO on a particular IC with the frequency of that RO on the same IC with the Trojan disabled. The computation is thus

$$TROI_{ROj,Ti} = (1/38)\sum_{k=1}^{k=38} \frac{|RO_{j,k,Tfree} - RO_{j,k,Ti}|}{RO_{j,k,Tfree}} \times 100\% \quad (7)$$

where $TROI_{ROj,Ti}$ is the mean impact of the $i^{th}$ Trojan on the $j^{th}$ RO across all ICs compared to the Trojan-free case. $RO_{j,k,Tfree}$ is the Trojan-free frequency for the $j^{th}$ RO on the $k^{th}$ IC, and similarly, $RO_{j,k,Tj}$ is the frequency of the $j^{th}$ RO on the $k^{th}$ IC with the $i^{th}$ Trojan activated.

It is with this calculation that the value of the single-IC multiple-Trojan design is best demonstrated. By comparing measurements made with a Trojan enabled against measurements made on the same IC with the Trojan disabled inter-die variation is eliminated from the analysis. Had separate ICs been fabricated with Trojans inserted and Trojans removed, only comparisons between different ICs would be possible and the computation would include inter-die process variation. By restricting comparisons to the same RO intra-die process variations are eliminated from the computation as well.

The results for Trojan impact are presented in Figure 5. It is immediately clear that Trojans of greater area and those which partially activate more frequently induce a greater change in the frequencies of nearby ROs since they consume more power. The maximum induced change for the largest Trojans in this experiment is representative of one of the core issues in the IC trust problem. The Trojan induces at most a change of 2.5% to frequencies, yet as Table 3 reports, intra-die variation and inter-die variation induce far greater changes suggesting these Trojans would be completely obfuscated in a test where these variations are not isolated. However, as discussed in Subsection 5.3, Trojan detection is still possible with this technique. The manner in which Trojan impact is distributed across ROs, including the decrease in impact on $RO_3$ and $RO_4$ for larger Trojans, is discussed in Subsection 5.2.

### 5.2 Spatial Locality Analysis

To analyze the effect of Trojan location, the ring oscillator which experiences the greatest Trojan impact calculated with Equation 7 is determined for each IC with a particular Trojan. A histogram

Figure 6: Number of instances of each RO being most impacted by a Trojan.

showing the frequency with which each ring oscillator was the most impacted on an IC is shown in Figure 6. The location of Trojan gates relative to the gates of the ROs and the vertical power line is shown in Figure 2

Notably, $RO_8$ is impacted most frequently for all Trojans since several of its gates are closest to the vertical power strap thereby causing a portion of the overall power supply noise to affect this RO. For $T_1$ and $T_2$ a substantial portion of the Trojan impact is distributed on $RO_2$ and $RO_3$ since these Trojans are located close to these ROs and likely share power lines.

Since the majority of the gates in subsequent Trojans are closest to $RO_8$, more of the Trojan impact is distributed on this RO. Perhaps counter-intuitively, the distribution becomes more focused on a single RO as the Trojan expands in size. Had the Trojan expanded vertically and towards multiple ROs it is likely the distribution would become less focused. However, for these Trojans which extend primarily horizontally, the increase in area and activity further increases the Trojan impact without expanding into other regions of the power network.

For $T_7$ the Trojan becomes less localized on $RO_8$ since $T_7$ is particularly close to the vertical power strap. For this reason, the Trojan impact is more evenly distributed across ROs since the vertical power strap supplies power to the entire circuit. Finally, the reduced impact on $RO_3$ and $RO_4$ for $T_6$ and $T_7$ shown in Figure 5 is due to the loosely distributed nature of these ROs away from the vertical power line and the placement of these Trojans close to the vertical power line.

## 5.3 IC Classification and False-Positive Analysis

In Section 5.1, it was shown that all Trojans used in this study impacted the RO frequencies substantially less than inter-die and intra-die process variations. However, using the principal component analysis (PCA) [16] based classification scheme presented below, it is still possible to detect these Trojans. In order to verify that this data is adequately represented in fewer than 8 principal components, the percent of the total variance in each PCA representation is computed by dividing the cumulative sum of the latent of the PCA representation by the total sum. The percent variance for each representation is shown in Table 4. The results imply that any representation of at least 2 components should adequately represent this data.

To succeed, a classification scheme must perform two functions: (1) it must correctly label Trojan-inserted circuits as tampered and (2) it must correctly label Trojan-free circuits as un-compromised. The steps for the presented classification scheme are:

Table 4: Percent variation contained in a representation of $h$ principal components.

| Components | Percent Variation |
| --- | --- |
| 1 | 89.4% |
| 2 | 99.39% |
| 3 | 99.59% |
| 4 | 99.79% |
| 5 | 99.87% |
| 6 | 99.93% |
| 7 | 99.97% |
| 8 | 100% |

1. Form a matrix from golden (Trojan-free) data in which each row is a verified Trojan-free IC and each column is a ring oscillator. Append a similar row containing the data from the chip under authentication (CUA) to the matrix.
2. Obtain a representation of this matrix using the first $h$ principal components
3. Render an $h$-dimensional convex hull [?] with all data except that of the CUA.
4. Determine if the CUA point falls within the hull. If it is within the boundaries of the hull it is considered Trojan-free.

To examine the performance of this classification scheme, the data are organized into five cases in which 8 of the 38 functioning ICs are randomly selected to represent Trojan-free chips to be authenticated and the remaining ICs are used to build the golden signature. All 38 ICs are used as Trojan-inserted chips under authentication.

The classification scheme was tested using both 2 and 3 dimensional hulls using the same subset cases for both hull types. The percent chips labeled as Trojan-inserted are shown for each case using both 2 and 3 dimensions are shown in Figure 7a and Figure 8a respectively. "FP" indicates the number of Trojan-free chips which were incorrectly classified. For both 2 and 3 dimensions, the behavior varies among the randomly selected cases. Thus for clarity, the average rates among all cases are shown in Figure 7b and Figure 8b. For both the 2 and 3 dimensional schemes, the false positive rates are lower than the detection rates for even the smallest Trojans in the experiment. For Trojans T1-T5 the detection rates are under 50%. This is unsurprising since these Trojans consisting of fewer than 130 transistors were intentionally designed to determine and emphasize the limitations of this technique.

For the larger Trojans, the detection rates are as high as 60-70% for the 2 dimensional case and 80-90% for the 3 dimensional case. Notably, the percent ICs labeled Trojan-inserted tends to be higher for the 3 dimensional case indicating sensitivity is related to the number of dimensions used. However, the three-dimensional case also achieves a higher ratio of detection rate to false positive rate for some cases.

These results demonstrate that the ring oscillator network scheme and the presented classification scheme can adequately separate Trojan-inserted designs from the Trojan-free designs despite the presence of obfuscating process variations. Although intra-die and inter-die variations introduce roughly 8% and 17% variations in RO frequencies respectively compared to the 1-3% change induced by the inserted Trojans, this technique successfully classifies ICs by exploiting the spatially correlated nature of process variations.

## 6. CONCLUSIONS

In this work, the RON structure for detecting hardware Trojans was analyzed using 38 ICs containing the ISCAS s9234 benchmark

(a) All cases using 2 dimensions



(b) Mean rates using 2 dimensions

Figure 7: Classification using the presented scheme and 2 dimensions.



(a) All cases using 3 dimensions



(b) Mean rates using 3 dimensions

Figure 8: Classification using the presented scheme and 3 dimensions.

circuit fabricated using the IBM 90nm process. We have shown that ring oscillator frequencies increase with increasing Trojan partial activity and that ring oscillators which share power lines with nearby Trojans will be most impacted. The presented results reveal that it is possible for Trojan impact to counter-intuitively become more localized as it expands in size provided it remains within the region most closely aligned with a single ring oscillator. Lastly, this work has demonstrated that even in the presence of obfuscating process variations, measurement noise, and environment variation ICs may still be effectively classified using a PCA-based classification technique. Future work will improve the classification procedure, and we will explore the potential for techniques which do not require a golden model.

# 7. ACKNOWLEDGEMENT

# 8. REFERENCES

[1] "Report of the Defense Science Board Task Force on High Performance Microchip Supply," Defense Science Board, US DoD, *http://www.acq.osd.mil/dsb/reports/2005-02-HPMSi_Report_Final.pdf*, Feb, 2005.

[2] X. Zhang and M. Tehranipoor,"RON: An On-chip Ring Oscillator Network for Hardware Trojan Detection," *in Proc. Design, Automation, and Test in Europe (DATE)*, pp. 1-6, 2011.

[3] R. Karri, J. Rajendran, K Rosenfeild, M. Tehranipoor "Trustworthy Hardware: Identifying and Classifying Hardware Trojans", *IEEE Design and Test of Computers*, pp. 39-46, 2010

[4] D. Agrawal, S. Baktir, D. Karakoyunlu, P. Rohatgi, and B. Sunar, "Trojan Detection using IC Fingerprinting," in *in Proc. IEEE Symposium on Security and Privacy (SP)*, pp. 296-310, 2007.

[5] R. Rad, J. Plusquellic, and M. Tehranipoor, "Sensitivity Analysis to Hardware Trojans using Power Supply Transient Signals," *IEEE Int. Symposium on Hardware-Oriented Security and Trust (HOST)*, pp. 3-7, June, 2008.

[6] M. Potkonjak et al., "Hardware Trojan Horse Detection Using Gate-Level Characterization," *in Proc. Design Automation Conf. (DAC)*, ACM Press, pp. 688-693, 2009.

[7] X. Wang, H. Salmani, M. Tehranipoor, and J. Plusquellic, "Hardware Trojan Detection and Isolation using Current Integration and Localized Current Analysis," in *in Proc. IEEE International Symposium on Defect and Fault Tolerance of VLSI Systems (DFT)*, pp. 87-95, 2008.

[8] Y. Jin and Y. Makris, "Hardware Trojan Detection using Path Delay Fingerprint," *in Proc. IEEE International Symposium on Hardware-Oriented Security and Trust (HOST)*, pp. 51-57, 2008.

[9] J. Li and J. Lach, "At-Speed Delay Characterization for IC Authentication and Trojan Horse Detection," *in Proc. IEEE Int. Hardware-Oriented Security and Trust (HOST)*, pp.8-14, 2008.

[10] M. Tehranipoor and F. Koushanfar, "A Survey of Hardware Trojan Taxonomy and Detection," *IEEE Design and Test of Computers*, pp. 10-25, 2010.

[11] S. Jha and S. K. Jha, "Randomization Based Probabilistic Approach to Detect Trojan Circuits," *in Proc. IEEE High Assurance System Engineering Symposium*, pp. 117-124, 2008.

[12] M. Banga and M. Hsiao, "A Region based Approach for the Identification of Hardware Trojans," *in Proc. IEEE Int.Symposium on Hardware-Oriented Security and Trust (HOST)*, pp. 40-47, 2008.

[13] F. Wolff, C. Papachristou, S. Bhunia, and R. S. Chakraborty, "Towards Trojan-free Trusted ICs: Problem Analysis and Detection Scheme" in *in Proc. Design, Automation and Test in Europe (DATE)*, pp. 1362-1365, 2008.

[14] H. Salmani, M. Tehranipoor, and J. Plusquellic, "A Novel Technique for Improving Hardware Trojan Detection and Reducing Trojan Activation Time,"

[15] M. Abramovici and P. Bradley, "Integrated Circuit Security: new Threats and Solutions," in *5th Annual Workshop on Cyber Security and information intelligence Research : Cyber Security and information intelligence Challenges and Strategies*, pp. 13-15, April. 2009.

[16] I. T. Jolliffe, "Principal Component Analysis (2ed Edition)," Springer, pp. 150-165, 2002.